

Ανάλυση Υπολοίπων

Τα υπόλοιπα για το πολλαπλό γραμμικό μοντέλο της παλινδρόμησης ορίζονται από τη σχέση $\underline{e} = \underline{Y} - \hat{\underline{Y}} = \underline{Y} - \underline{X} \cdot \hat{\underline{\beta}} = (\underline{I} - \underline{P}) \cdot \underline{Y}$, όπου $\underline{P} = \underline{X} \cdot (\underline{X}^T \cdot \underline{X})^{-1} \cdot \underline{X}^T$.

Ιδιότητες των υπολοίπων:

$$(1) \quad E(\underline{e}) = \underline{0}$$

(απόδειξη)

$$E(\underline{e}) = E[(\underline{I} - \underline{P}) \cdot \underline{Y}] = (\underline{I} - \underline{P}) \cdot E(\underline{Y}) = (\underline{I} - \underline{P}) \cdot \underline{X} \cdot \underline{\beta} = (\underline{I} \cdot \underline{X} - \underline{P} \cdot \underline{X}) \cdot \underline{\beta} = \underline{0} \cdot \underline{\beta}$$

$$(2) \quad Cov(\underline{e}) = \sigma^2 \cdot (\underline{I} - \underline{P})$$

(απόδειξη)

$$\begin{aligned} Cov(\underline{e}) &= Cov[(\underline{I} - \underline{P}) \cdot \underline{Y}] = (\underline{I} - \underline{P}) \cdot Cov(\underline{Y}) \cdot (\underline{I} - \underline{P})^T = (\underline{I} - \underline{P}) \cdot \sigma^2 \cdot \underline{I} \cdot (\underline{I} - \underline{P})^T = \\ &= \sigma^2 \cdot (\underline{I} - \underline{P}) \cdot (\underline{I} - \underline{P})^T = \sigma^2 \cdot (\underline{I} - \underline{P}) \cdot (\underline{I} - \underline{P}) = \sigma^2 \cdot (\underline{I} - \underline{P})^2 = \sigma^2 \cdot (\underline{I} - \underline{P}) \end{aligned}$$

Άρα συμπεραίνουμε ότι $Var(e_i) = \sigma^2 \cdot (1 - P_{ii})$ και $Cov(e_i, e_j) = -\sigma^2 \cdot P_{ij}$, $i \neq j$, όπου $\underline{P} = (P_{ij})$, $i = 1, 2, \dots, n$, $j = 1, 2, \dots, n$.

(3) Κατανομή των υπολοίπων:

$$\underline{e} \sim N_n(\underline{0}, \sigma^2 \cdot (\underline{I} - \underline{P}))$$

(απόδειξη)

Αφού $\underline{e} = (\underline{I} - \underline{P}) \cdot \underline{Y}$ είναι γραμμικός συνδυασμός κανονικών.

$$(4) \quad Cov(\underline{e}, \hat{\underline{Y}}) = \underline{0}$$

(απόδειξη)

$$\begin{aligned} Cov(\underline{e}, \hat{\underline{Y}}) &= Cov((\underline{I} - \underline{P}) \cdot \underline{Y}, \underline{P} \cdot \underline{Y}) = (\underline{I} - \underline{P}) \cdot Cov(\underline{Y}, \underline{Y}) \cdot \underline{P} = (\underline{I} - \underline{P}) \cdot Cov(\underline{Y}) \cdot \underline{P} = \\ &= (\underline{I} - \underline{P}) \cdot \sigma^2 \cdot \underline{I} \cdot \underline{P} = \sigma^2 \cdot (\underline{I} - \underline{P}) \cdot \underline{P} = \sigma^2 \cdot (\underline{P} - \underline{P}^2) = \sigma^2 \cdot (\underline{P} - \underline{P}) = \sigma^2 \cdot \underline{0} = \underline{0} \end{aligned}$$

Άρα τα υπόλοιπα και οι προβλέψεις είναι ασυσχέπιστα και όταν ισχύει η κανονικότητα, είναι και ανεξάρτητες τυχαίες μεταβλητές.

$$(5) \quad Cov(\underline{e}, \underline{Y}) = \sigma^2 \cdot (\underline{I} - \underline{P})$$

(απόδειξη)

$$\begin{aligned} Cov(\underline{e}, \underline{Y}) &= Cov((\underline{I} - \underline{P}) \cdot \underline{Y}, \underline{Y}) = (\underline{I} - \underline{P}) \cdot Cov(\underline{Y}, \underline{Y}) = (\underline{I} - \underline{P}) \cdot Cov(\underline{Y}) = (\underline{I} - \underline{P}) \cdot \sigma^2 \cdot \underline{I} = \\ &= \sigma^2 \cdot (\underline{I} - \underline{P}) \end{aligned}$$

Δύο άλλες μορφές υπολοίπων:

(α) Τυποποιημένα Υπόλοιπα (standardized residuals):

$$e_{si} = \frac{e_i}{\hat{\sigma}} = \frac{Y_i - \hat{Y}_i}{\hat{\sigma}}$$

Τα τυποποιημένα υπόλοιπα δεν εξαρτώνται από τη μονάδα μέτρησης και δείχνουν το μέγεθος των e_i σε σχέση με $\hat{\sigma}$.

(β) Μαθηματικοποιημένα Υπόλοιπα (studentized residuals):

$$e_{sti} = \frac{e_i}{\hat{\sigma} \cdot \sqrt{1 - P_{ii}}}$$

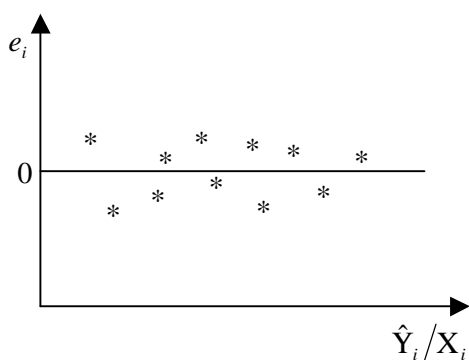
Ανάλυση Υπολοίπων

Η ανάλυση υπολοίπων είναι μια σειρά μεθόδων (γραφικών και στατιστικών) που μας επιτρέπει να ελέγξουμε τις υποθέσεις για τα σφάλματα, καθώς και την ορθότητα ή μη του μοντέλου.

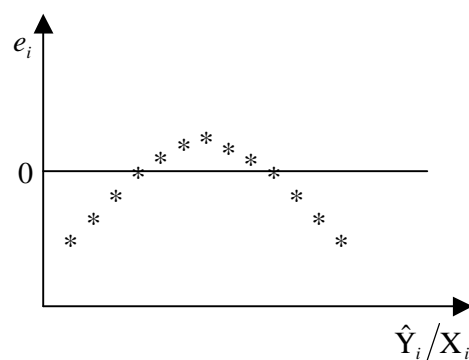
A) Έλεγχος Ύπαρξης Ασυσχέτιστων Σφαλμάτων:

Μια γραφική μέθοδος είναι η εξής:

Κατασκευάζουμε γραφική παράσταση των υπολοίπων e_i ως προς \hat{Y}_i ή/και X_i , επειδή τα σφάλματα είναι ασυσχέτιστα με \hat{Y}_i και X_i . Δεν κατασκευάζω γραφική παράσταση των e_i ως προς Y_i , γιατί δεν είναι ασυσχέτιστα.



(1)



(2)

Στο (2) έχω ένδειξη συσχισμένων σφαλμάτων, ενώ στο (1) έχω ένδειξη ασυσχέτιστων σφαλμάτων.

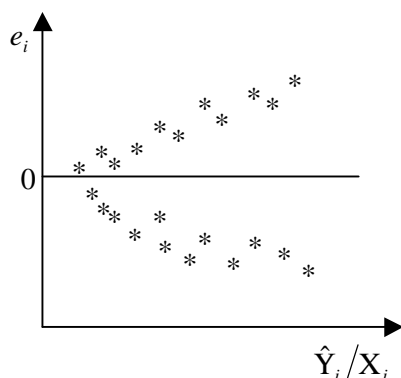
Ένας στατιστικός έλεγχος είναι ο παρακάτω: (έλεγχος των ροών)

H_0 : η ακολουθία των "+" στα υπόλοιπα είναι τυχαία (ασυσχέτιστα σφάλματα)

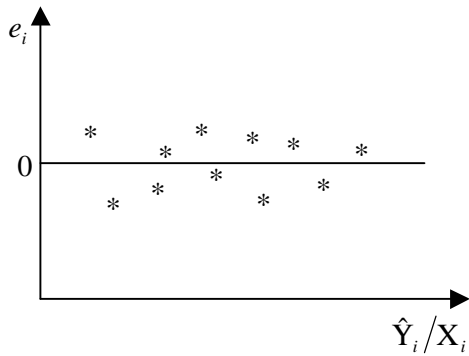
H_1 : η ακολουθία των "+" στα υπόλοιπα δεν είναι τυχαία

B) Έλεγχος Σταθερής Διακύμανσης: ($Var(\varepsilon_i) = \sigma^2, i = 1, 2, \dots, n$)

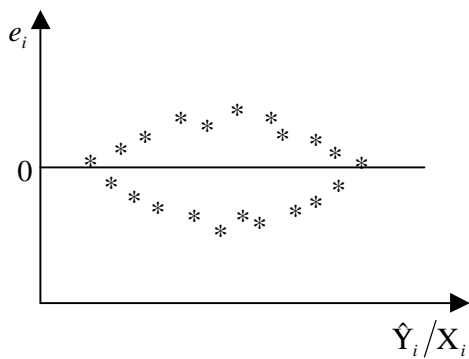
Σταθερή διακύμανση συνεπάγεται ότι η διακύμανση των σφαλμάτων δεν εξαρτάται από τα X_i (τιμές της ανεξάρτητης μεταβλητής).



Για μεγάλες τιμές των X_i έχω μεγάλες αποκλίσεις των σφαλμάτων. Ένδειξη μη σταθερής διακύμανσης.



Εδώ βλέπουμε σταθερή διακύμανση.



Για μεσαίες τιμές των X_i έχω μεγάλες αποκλίσεις. Για ακραίες τιμές έχω μικρές αποκλίσεις. Ένδειξη μη σταθερής διακύμανσης.

Διόρθωση μη σταθερής διακύμανσης:

Υπάρχουν οι εξής δύο μέθοδοι:

- (1) Μετασχηματισμός σταθεροποίησης της διακύμανσης.
- (2) Γενικευμένοι εκτιμητές ελαχίστων τετραγώνων.

Μέθοδος #1:

Έστω το απλό γραμμικό μοντέλο $Y_i = \beta_0 + \beta_1 \cdot X_i + \varepsilon_i$, $i = 1, 2, \dots, n$, με $E(\varepsilon_i) = 0$, $Cov(\varepsilon_i, \varepsilon_j) = 0$, $i \neq j$, και $Var(\varepsilon_i) = \sigma^2 \cdot X_i$ (εξαρτάται από X_i).

Χρησιμοποιούμε τις εκτιμήτριες ελαχίστων τετραγώνων $\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \cdot \bar{X}$ και

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (X_i - \bar{X}) \cdot (Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2}, \text{ αλλά οι ιδιότητες τους δεν ισχύουν πλέον } (\hat{\beta}_0, \hat{\beta}_1 \text{ δεν}$$

είναι ΑΟΕΔ όταν ισχύει η κανονικότητα).

Εφαρμόζω κάποιο μετασχηματισμό για να βρώ άλλο εκτιμητή:

$$\text{Αν όλα τα } X_i \neq 0, \frac{Y_i}{X_i} = \frac{\beta_0}{X_i} + \beta_1 + \frac{\varepsilon_i}{X_i} \Rightarrow Y_i^* = \beta_0^* + \beta_1^* \cdot X_i^* + \varepsilon_i^*,$$

$$\text{όπου } Y_i^* = \frac{Y_i}{X_i}, X_i^* = \frac{1}{X_i}, \beta_0^* = \beta_1, \beta_1^* = \beta_0, \varepsilon_i^* = \frac{\varepsilon_i}{X_i}.$$

Αλλά τώρα έχω το μοντέλο $Y_i^* = \beta_0^* + \beta_1^* \cdot X_i^* + \varepsilon_i^*$, με

$$\left. \begin{array}{l} E(\varepsilon_i^*) = 0 \\ \text{Var}(\varepsilon_i^*) = \sigma^2 \\ \text{Cov}(\varepsilon_i^*, \varepsilon_j^*) = 0, i \neq j \\ \varepsilon_i^* \sim N(0, \sigma^2) \end{array} \right\} \begin{array}{l} \text{κλασσικές υποθέσεις} \\ \text{του γραμμικού μοντέλου} \end{array}$$

Επομένως, οι εκτιμήτριες ελαχίστων τετραγώνων $\hat{\beta}_0^*, \hat{\beta}_1^*$ έχουν όλες τις βέλτιστες ιδιότητες και συνεπώς $\hat{\beta}_1 = \hat{\beta}_0^*$ και $\hat{\beta}_0 = \hat{\beta}_1^*$.

Μέθοδος #2:

Γενικά μπορούμε να θεωρήσουμε το πολλαπλό μοντέλο παλινδρόμησης $\underline{Y} = \underline{X} \cdot \underline{\beta} + \underline{\varepsilon}$, $E(\underline{\varepsilon}) = \underline{0}$, $\text{Cov}(\underline{\varepsilon}) = \sigma^2 \cdot \underline{\Sigma}$, $\underline{\Sigma} > 0$, όπου $\underline{\Sigma}$ γνωστός $(n \times n)$ πίνακας. Άρα τώρα έχω συσχετισμένα σφάλματα.

$$\text{Αφού } \underline{\Sigma} > 0 \Rightarrow \exists \underline{\Sigma}^{-1} > 0 \Rightarrow \exists \underline{\Sigma}^{-1/2} : \underline{\Sigma}^{-1/2} \cdot \underline{\Sigma}^{-1/2} = \underline{\Sigma}^{-1}.$$

$$\text{Οπότε } \underline{\Sigma}^{-1/2} \cdot \underline{Y} = \underline{\Sigma}^{-1/2} \cdot \underline{X} \cdot \underline{\beta} + \underline{\Sigma}^{-1/2} \cdot \underline{\varepsilon} \Rightarrow \underline{Y}^* = \underline{X}^* \cdot \underline{\beta} + \underline{\varepsilon}^*,$$

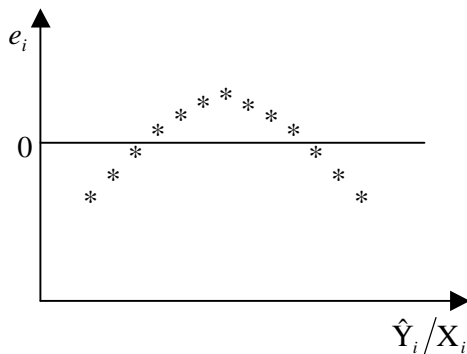
όπου $\underline{\varepsilon}^* \sim N_n(\underline{0}, \sigma^2 \cdot \underline{I})$, όταν ισχύει η κανονικότητα.

$$\begin{aligned} \text{Συνεπώς, } \hat{\underline{\beta}} &= (\underline{X}^{*\text{T}} \cdot \underline{X}^*)^{-1} \cdot \underline{X}^{*\text{T}} \cdot \underline{Y}^* = (\underline{X}^{\text{T}} \cdot \underline{\Sigma}^{-1/2} \cdot \underline{\Sigma}^{-1/2} \cdot \underline{X})^{-1} \cdot \underline{X}^{\text{T}} \cdot \underline{\Sigma}^{-1/2} \cdot \underline{\Sigma}^{-1/2} \cdot \underline{Y} = \\ &= (\underline{X}^{\text{T}} \cdot \underline{\Sigma}^{-1} \cdot \underline{X})^{-1} \cdot \underline{X}^{\text{T}} \cdot \underline{\Sigma}^{-1} \cdot \underline{Y} \quad (\text{γενικευμένη εκτιμήτρια ελαχίστων τετραγώνων}) \end{aligned}$$

Η γενικευμένη εκτιμήτρια ελαχίστων τετραγώνων έχει όλες τις βέλτιστες ιδιότητες, π.χ. $\hat{\underline{\beta}}$ είναι ΑΟΕΔ.

Γ) Έλεγχος Ορθότητας Γραμμικού Μοντέλου:

Γραμμική μέθοδος: Γραφική παράσταση των υπολοίπων ως προς X_i ή ως προς \hat{Y}_i .



Παρατηρείται ο σχηματισμός κάποιας καμπύλης. Πιθανώς να πρέπει να εισαχθεί κάποιος δευτεροβάθμιος όρος ή τα σφάλματα να είναι συσχετισμένα.

Δηλαδή: Ενώ υποθέτουμε ότι ισχύει το γραμμικό μοντέλο $Y = \beta_0 + \beta_1 \cdot X + \varepsilon$, το πιο πάνω σχήμα δίνει ένδειξη του μοντέλου $Y = \beta_0 + \beta_1 \cdot X + \beta_2 \cdot X^2 + \varepsilon$.

Εξήγηση: Αφού $Y - \hat{Y} = (\beta_0 - \hat{\beta}_0) + (\beta_1 - \hat{\beta}_1) \cdot X$, δεν θα πρέπει να έχω καμία συστηματική εμφάνιση στο σχήμα. Αλλά αν ισχύει το δευτεροβάθμιο μοντέλο, $Y - \hat{Y} = (\beta_0 - \hat{\beta}_0) + (\beta_1 - \hat{\beta}_1) \cdot X + \beta_2 \cdot X^2$.

Παρατήρηση:

Συχνά, ένα μη γραμμικό μοντέλο μπορεί να μετασχηματιστεί σε γραμμικό.

π.χ. $Y = \alpha \cdot e^{\beta \cdot X} \Rightarrow \log(Y) = \log(\alpha) + \beta \cdot X$
 $Y = \alpha \cdot X^p \Rightarrow \log(Y) = \log(\alpha) + p \cdot \log(X)$

Δ) Έλεγχος Κανονικότητας Σφαλμάτων:

Αυτός ο έλεγχος εκτελείται με έλεγχο κανονικότητας των υπολοίπων ή με γραφική παράσταση (Q-Q plot).

Αν η τ.μ. X είναι κανονική, τότε $P(X \leq x_p) = 1 - p \Rightarrow P\left(Z \leq \frac{x_p - \mu}{\sigma}\right) = 1 - p$.

Άρα $\frac{x_p - \mu}{\sigma} = z_p \Rightarrow x_p = \mu + \sigma \cdot z_p$.

Συνεπώς,



Θα έχω ευθεία γραμμή αν το δείγμα προέρχεται από την κανονική κατανομή.

Παρατήρηση:

Αν κάνουμε γραφικές παραστάσεις των e_{si} ή e_{sti} , καταλήγουμε στα ίδια συμπεράσματα.

Ακραίες Παρατηρήσεις

Η παρατήρηση Y_j (j συγκεκριμένο) είναι ακραία αν το αντίστοιχο υπόλοιπο $e_j = Y_j - \hat{Y}_j$ είναι κατά πολύ μεγαλύτερο (κατά απόλυτη τιμή) από όλα τα άλλα υπόλοιπα.

Μεθοδολογία για την ανεύρεση ακραίας παρατήρησης:

A) Παραλείπουμε από την ανάλυση την j παρατήρηση. Έτσι έχουμε $(n-1)$ παρατηρήσεις.

B) Σχηματίζουμε τις εκτιμήτριες $\hat{\beta}_{(-j)}, \hat{\sigma}_{(-j)}^2$, που στηρίζονται στις $(n-1)$ παρατηρήσεις.

Γ) Σχηματίζουμε τη διαφορά $Y_j - \hat{Y}_{(-j)}$, όπου η $\hat{Y}_{(-j)}$ προκύπτει από τις $(n-1)$ παρατηρήσεις, δηλαδή $\hat{Y}_{(-j)} = \underline{X}_j^T \cdot \hat{\beta}_{(-j)}$, όπου \underline{X}_j^T η j γραμμή του πίνακα \mathbf{X} .

Δ) Ελέγχω αν η j παρατήρηση είναι ακραία με ελεγχοσυνάρτηση

$$t = \frac{Y_j - \hat{Y}_{(-j)}}{\hat{\sigma}_{(-j)} \cdot \sqrt{A}} \sim t_{(n-1)-(p+1)}, \text{ όπου } A = \left[1 + \underline{X}_j^T \cdot (\mathbf{X}_{(-j)}^T \cdot \mathbf{X}_{(-j)})^{-1} \cdot \underline{X}_j \right].$$

Απορρίπτουμε την υπόθεση H_0 : η j παρατήρηση δεν είναι ακραία, σε επίπεδο σημαντικότητας α , αν $|t| > t_{(n-1)-(p+1), \alpha/2}$.